Any general code comments, questions, tips/tricks?

# R can play sound, text-to-talk, or open a youtube video?

```
# package beepr can choose from ~10 different sound bytes
beepr::beep(sound = "fanfare")

# system just calls to command line
system("say Worms are cool!")

# you can also set to open a youtube video (or any link)
system("open https://www.youtube.com/watch?v=xyNICGyYReQ")
# or just > open https://www.youtube.com/watch?v=xyNICGyYReQ on terminal
```

*Or just download youtube videos with terminal: here*

# Git/Github authentication setup?

**Set credentials on Mac:**

```
git config --global credential.helper osxkeychain
```

**Erase credentials on Mac:**

```
git credential-osxkeychain erase
host=github.com
protocol=https
```

*Need a personal authentication token: here*

# Git/Github authentication setup?

**Set credentials on QUEST:**

```
git config --global credential.helper cache
```

**Erase credentials on QUEST:**

```
git config --global --unset credential.helper
```

*Need a personal authentication token: here*

# Other misc. finds

```r
# filter any column for "zinc"
zincdf <- gene_descriptions %>%
    dplyr::filter_all(any_vars(str_detect(., pattern = "zinc")))

# extract gene ids!!!
stringr::str_extract(info, "WBGene........")

# remove last comma from string
gsub(',$', '', test)
```

*Maybe next time we should talk about* `stringr`*?*

# Running nextflow pipelines with GitHub

1. Clone git repo to QUEST

   ```
   git clone https://github.com/AndersenLab/NemaScan.git
   cd NemaScan
   ```

2. Run pipeline

   ```
   nextflow run develop.nf --debug
   ```

Local

1. Run pipeline remotely using git repo name

   ```
   nextflow run AndersenLab/NemaScan/develop.nf --debug
   ```

Remote

# Running nextflow pipelines with GitHub

PRO

CON

Local

- Only update pipeline when you manually `git pull` (control)
- Easy to make changes to suit your specific needs

- Might end up with several different versions of the same pipeline
- Might not know if there is a new update

Remote

- Always get newest version (but can still run older versions)
- Nextflow logs git commit - reproducible versions
- Can run anywhere on QUEST

- Will still need to clone/branch to make specific edits

# How to run a specific git branch/commit with nextflow

\# run master branch, newest commit

```
nextflow run AndersenLab/NemaScan --debug
```

\# run specific commit (maybe a previous version)

```
nextflow run andersenlab/nemascan --debug -r <commitid>
```

\# run newest commit on specific branch

```
nextflow run AndersenLab/NemaScan --debug -r cendr_dev
```
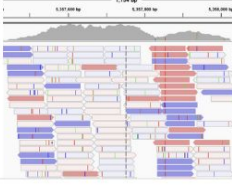
➔ *More info: [here](#)*

# Why would I want to run nextflow remote?



```
log.txt                                                          Open with TextEdit

Pipeline execution summary
---------------------------
Completed at: 2021-08-12T12:54:58.916823-05:00
Duration    : 4m 48s
Success     : false
workDir     : /projects/b1042/AndersenLab/work
exit status : null
Error report: SIGINT
Git info: https://github.com/AndersenLab/NemaScan.git - master [1fcb670fd4cb752532353e8f0eb3917cb4e60ab3]

{ Parameters }
---------------------------
Phenotype File                           = /home/kek973/.nextflow/assets/AndersenLab/NemaScan/input_data/elegans/phenotypes/abamectin_pheno.tsv
VCF                                      = 330_TEST.vcf.gz
Significance Threshold                   = BF
P3D                                      = TRUE
Max AF for Burden Mapping                = 0.05
Min Strains with Variant for Burden      = 2
Threshold for grouping QTL               = 1000
Number of SNVs to define CI              = 150
Eigen Memory allocation                  = 10 GB
Path to R libraries.                     = /projects/b1059/software/R_lib_3.6.0
Mapping                                  = RUN
Simulation                               = null
Simulate QTL effects                     = null
Annotation                               = null
Result Directory                         = Analysis_Results-20210812
```
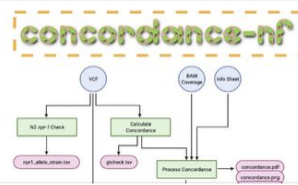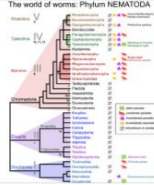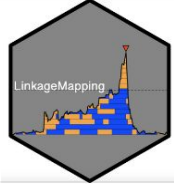
# What pipelines *should* I run remote?
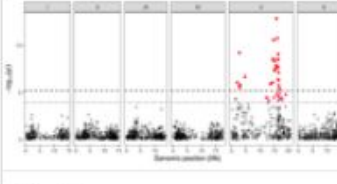

Alignment-nf


cegwas2-nf


concordance-nf


genomes-nf


linkagemapping-nf (a...


nil-ril-nf


Nemascan


post-gatk-nf


trim-fq-nf


wi-gatk

# Passing filename without file extension as SBATCH variable

```
# You have a set of files with the same file extension
# (e.g.: VX34.fa)
# You want to run a script on each one of this files

for file in *.fa; do sbatch --export=file=$file script.sh; done
```

```
# Let's say that you have pairs of files with same file name but different extensions
# (e.g.: VX34.fa, VX34.gff)
# You can remove the last file extension using ${var%.*} within this command:

for file in *.fa; do sbatch --export=file=${file%.*} script.sh; done

# ${var%.*} removes the shortest match to ".*" from $var
# You can use any extension delimiter (e.g: ${var%-*})
# You can then reference both files inside the script as $file.fa and $file.gff
```

*Credit: Nic*

# Passing filename without file extension as SBATCH variable

```
# Let's say you have pairs of files with same file name but with multiple different extensions
(e.g: VX34.filtered_contigs.m40.fa, VX34.curated.only.gff3)
# You can remove all file extensions using ${var%%.*} within this command:

for file in *.fa; do sbatch --export=file=$file,idef=${file%%.*} script.sh; done

# ${var%%.*} removes the longest match to ".*" from $var
# You can use any extension delimiter (e.g: ${var%%-*})
# You can then reference both files inside the script as $file and $idef.curated.only.gff
```

# Using the previous expression to generate lists

```
# Let's say you have a folder full of FASTQ, and you want to generate a strain list from that folder.
# Most FASTQ files in our storage are formatted as 'STRAIN_[...].fq.gz'
# -(e.g: XZ1733_CKDL200150306-1a-GA03-AK1863_HV53CDSXX_L3_2P.fq.gz)
# You can write a strain list using the ${var%%_*} within this command:

for file in *.fq.gz; do echo ${file%%_*}; done > strain-list.txt
```

```
# If you are expecting duplicates (different lanes for the same strain) you can pipe 'sort -u':

for file in *.fq.gz; do echo ${file%%_*}; done | sort -u > strain-list.txt
```

*Credit: Nic*

# Passing variable from list as SBATCH variable

```
# You got your strain list, and wish to run a script for each line in your list
# You can loop through the file and pass variables using the following command

while IFS= read -r strain; do sbatch --export=strain=$strain script.sh; done < strain-list.txt
```

# Generate strain list from from directory list

```
# Let's say Katie sent you a list of full paths to FASTQs of interest,
# and you wish to simplify it into a strain list...
# You can loop through the file, extract the file names, and extract the
# strain name using the following command

while IFS= read -r path;
do filename=$(basename $path); strain=${filename%%_*}; echo $strain;
done < fastq-directories.txt | sort -u > strain-list.txt
```

*Credit: Nic*

# Subsetting with grep (different file types)

```
# Let's say you want to subset a file based on whole-word matches from a list
# -e.g.: subsetting your sample sheet based on a list of strainst

grep -wFf subset.txt all.txt > all_subset.txt
```

```
# or if applied to the example above:

grep -wFf strains.txt sample_sheet.tsv > subset_ss.tsv

# This will return every line in 'sample_sheet.tsv' that has a whole-word match
# (literal substring) to each line in 'strains.txt'
# Doesn't require sorted files, and lines can have different fields
```
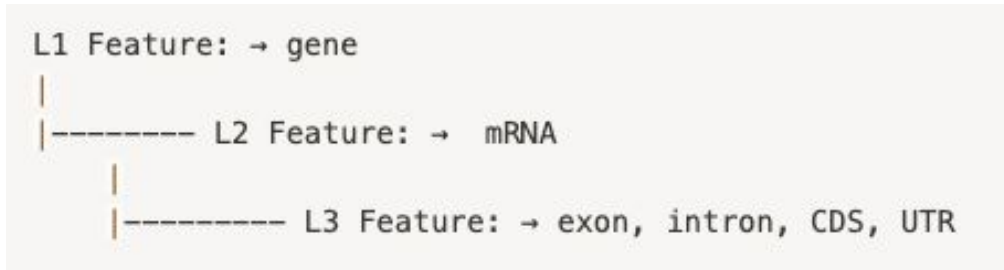
*Credit: Nic*

# Subsetting & pattern matching with comm (same file types)

```
# uses whole-line matches and files must be sorted

comm -12 file1 file2 > common
comm -23 file1 file2 > file1_uniques
comm -13 file1 file2 > file2_uniques
```

*Credit: Nic*

# Manipulating wormbase GFFs with grep & sed

- GFF annotations are flat files that have the following structure for protein coding genes:

```
L1 Feature: → gene
 |
 |--------- L2 Feature: →  mRNA
      |
      |--------- L3 Feature: → exon, intron, CDS, UTR
```

- Low level features are nested within higher level features, with genes being at the top
- Genes have different 'Attributes' than lower level features. The Attributes field points to the parent in lower level features, and provides additional information about that gene in L1 features.

*Credit: Nic*

# Manipulating wormbase GFFs with grep & sed

- Sample gene in GFF:



- If you tried to use grep to pull the GFF features using a wormbase ID for a given gene, it will only return the L1/L2 features, but none of the L3 features



*Credit: Nic*

# Manipulating wormbase GFFs with grep & sed

- A simple solution would be to find a common attribute across all features (like gene ID), and pull all lines pertaining to that gene

```
grep "Y74C9A.1" annotation.gff > gene.gff
```

- This would work great for a single gene, but what if you had a massive list of WB genes? Too much manual work…

# Manipulating wormbase GFFs with grep & sed

```
# Let's say you have a list of WB genes called "WBids.txt", you can use this one-liner:
grep -wFf WBids.txt WB.annotation.gff3 | \   # this returns L1/L2 features that match gene IDs
grep -o "sequence_name.*" | \ # Returns only L1 features, omitting every field prior to the gene ID.
sed 's/\;.*//g' | \ #This removes any fields after the first field in each line, delimited by ';'
sed 's/sequence_name=//g' > gene_ids.txt #This cleans up the attribute prefix, leaving the gene ID.

#You can then filter the original GFF with your new transcript list!

grep -wFf gene_ids.txt WB.annotation.gff3 > genes.gff3
```

*Credit: Nic*

# Manipulating wormbase GFFs with grep & sed

```
# Another implementation of this code is filtering by biotype
# Wormbase GFFs contain multiple L1 features aside from genes (e.g.: ncRNA)
# Protein coding genes have the "biotype=protein_coding"
# Let's say that you would like to filter out any non protein coding features:

grep "protein_coding" WB.annotation.gff3 | \  # this returns L1 features that match the biotype
grep -o "sequence_name.*" | \ # Returns only L1 feature attributes, omitting every field prior to th
sed 's/\;.*//g' | \ #This removes any fields after the first field in each line, delimited by ';'
sed 's/sequence_name=//g' > gene_ids.txt #This cleans up the attribute prefix, leaving the gene ID.

grep -wFf gene_ids.txt WB.annotation.gff3 > protein_coding.gff3
```

*Credit: Nic*